

TIA 2024

Predictia finală (cea mai bună): 0.85810

Metoda de învățare folosită

Am ales metoda **Random Forest**. Aceasta se bazează pe crearea și combinarea mai multor arbori de decizie pentru a îmbunătăți acuratețea și robustețea modelului.

Ne imaginăm o pădure formată din mai mulți arbori. Fiecare arbore reprezintă o opinie independentă despre datele acestea. În final, Random Forest ia decizia pe baza majorității acestor opinii, ceea ce duce la o predicție robustă și precisă.

Astfel, Random Forest combină puterea mai multor modele slabe (arbori de decizie) pentru a crea un model puternic și stabil.

Codul

```
1  import numpy as np
2  import pandas as pd
3  from sklearn.model_selection import train_test_split
4  from sklearn.preprocessing import StandardScaler
5  from sklearn.decomposition import PCA
6  from sklearn.ensemble import RandomForestClassifier
7  from sklearn.metrics import precision_score
8
9  class FashionMNISTAnalyzer:
10     def __init__(self, train_path, test_path, submission_path):
11
12         train_data = np.load(train_path)
13         self.x_train = train_data['x_train']
14         self.y_train = train_data['y_train']
15
16         test_data = np.load(test_path)
17         self.x_test = test_data["x_test"]
18
19         self.submission = pd.read_csv(submission_path)
20
21     def preprocess_data(self, subset_ratio=0.8, apply_pca=True, n_components=100):
22
23         x_train, _, y_train, _ = train_test_split(
24             self.x_train, self.y_train,
25             train_size=subset_ratio,
26             random_state=42
27         )
28
29         x_train, x_val, y_train, y_val = train_test_split(
30             x_train, y_train,
31             test_size=0.2,
32             random_state=42
33         )
34
35         scaler = StandardScaler()
36         x_train_scaled = scaler.fit_transform(x_train)
37         x_val_scaled = scaler.transform(x_val)
38         x_test_scaled = scaler.transform(self.x_test)
39
40         if apply_pca:
41             pca = PCA(n_components=n_components, random_state=42)
42             x_train_scaled = pca.fit_transform(x_train_scaled)
43             x_val_scaled = pca.transform(x_val_scaled)
44             x_test_scaled = pca.transform(x_test_scaled)
45
46         return (
47             x_train_scaled, x_val_scaled,
48             y_train, y_val,
49             x_test_scaled
50         )
51
52     def train_random_forest(self, x_train, y_train):
53
54         rf_classifier = RandomForestClassifier(
55             n_estimators=200,
56             max_depth=25,
```

```
56         max_depth=25,
57         min_samples_split=2,
58         min_samples_leaf=1,
59         bootstrap=True,
60         random_state=42,
61         n_jobs=-1
62     )
63     rf_classifier.fit(x_train, y_train)
64     return rf_classifier
65
66 def evaluate_model(self, model, x_val, y_val):
67
68     y_pred = model.predict(x_val)
69     precision = precision_score(y_val, y_pred, average='macro')
70     print(f"Precision: {precision:.4f}")
71     return precision
72
73 def generate_predictions(self, model, x_test):
74
75     predictions = model.predict(x_test)
76     self.submission['Label'] = predictions
77     print("Updated sample_submission.csv with predictions.")
78     return self.submission
79
80 def main():
81
82     train_path = "train.npz"
83     test_path = "test.npz"
84     submission_path = "sample_submission.csv"
85
86     analyzer = FashionMNISTAnalyzer(train_path, test_path, submission_path)
87
88     x_train, x_val, y_train, y_val, x_test = analyzer.preprocess_data(
89         subset_ratio=0.8,
90         apply_pca=True,
91         n_components=100
92     )
93
94     print("Training Random Forest:")
95     rf_model = analyzer.train_random_forest(x_train, y_train)
96
97     print("Evaluating Random Forest Model:")
98     precision = analyzer.evaluate_model(rf_model, x_val, y_val)
99     print(f"Final Precision Score: {precision:.4f}")
100
101    submission = analyzer.generate_predictions(rf_model, x_test)
102    submission.to_csv('final_submission.csv', index=False)
103
104 if __name__ == "__main__":
105     main()
```

Explicarea Codului

Importarea bibliotecilor necesare:

Clasa FashionMNISTAnalyzer:

- Metoda `init`: Încarcă datele de antrenament și testare din fișiere .npz și şablonul de trimitere dintr-un fișier .csv
- Metoda `preprocess_data`: Preprocesează datele prin împărțirea lor în seturi de antrenament și validare, scalarea caracteristicilor și aplicarea PCA pentru reducerea dimensionalității
- Metoda `train_random_forest`: Antrenează un model Random Forest folosind RandomizedSearchCV pentru a găsi cei mai buni hiperparametri
- Metoda `evaluate_model`: Evaluează modelul antrenat pe setul de validare și calculează scorul de precizie
- Metoda `generate_predictions`: Generează predicții pentru setul de testare și actualizează fișierul de trimitere cu aceste predicții

Funcția main:

- Definește căile către fișierele de date
- Inițializează un obiect FashionMNISTAnalyzer
- Preprocesează datele
- Antrenează modelul Random Forest
- Evaluează modelul și afișează scorul de precizie
- Generează predicții pentru setul de testare și salvează rezultatele într-un fișier .csv

Explicația folosirii Random Forest-ului:

Random Forest este un algoritm de învățare automată utilizat pentru clasificare și regresie. Este compus dintr-o colecție de arbori de decizie antrenați pe diferite subseturi ale datelor de antrenament și utilizează medierea pentru a îmbunătăți acuratețea și a controla supraînvățarea (așa cum am explicitat și stabilit mai devreme).

Pașii principali în utilizarea Random Forest-ului în acest cod sunt:

- Definirea hiperparametrilor:
 - n_estimators: Numărul de arbori în pădurea aleatorie.
 - max_depth: Adâncimea maximă a fiecărui arbore.
 - min_samples_split: Numărul minim de eșantioane necesare pentru a împărți un nod.
 - min_samples_leaf: Numărul minim de eșantioane necesare pentru a fi la un nod frunză
- Antrenarea modelului:

Se utilizează RandomizedSearchCV pentru a efectua o căutare aleatorie a celor mai buni hiperparametri. Aceasta este o metodă eficientă din punct de vedere al timpului, deoarece nu testează toate combinațiile posibile de hiperparametri, ci doar un subset aleatoriu. Modelul este antrenat pe datele de antrenament preprocesate.

- Evaluarea modelului:

Modelul antrenat este evaluat pe setul de validare pentru a calcula scorul de precizie. Aceasta ajută la determinarea performanței modelului pe date nevăzute.

De ce precizia este 0.85810 și nu mai mare?

- Reducerea dimensiunii cu PCA:

Utilizarea PCA reduce dimensiunea datelor la 100 de componente principale. În timp ce PCA ajută la scăderea timpului de rulare și la eliminarea redundanței, poate duce la

pierderea unor informații relevante, mai ales în cazul imaginilor din Fashion MNIST, unde detaliile fine contează.

- Parametrii Random Forest:

Modelul utilizează 200 de arbori și o adâncime maximă de 25. Deși acești parametri sunt solizi, pot să nu fie suficient de puternici pentru a captura toate caracteristicile setului de date.

- Proporția de date folosită pentru antrenament:

Se utilizează doar 80% din date pentru antrenament și 20% pentru validare. Reducerea setului de antrenament poate diminua performanța generală, deoarece Random Forest beneficiază de volume mari de date.

- Scalarea cu StandardScaler:

În unele cazuri, imaginile în format brut (neprocesate prin PCA) pot beneficia mai mult de scalare fără reducerea dimensiunii.

Așadar, precizia este de 0.85810, însă putem modifica programul astfel încât să obținem o precizie mai bună prin optimizarea parametrilor, sau monitorizarea timpului de rulare, sau creșterea numărului de componente, sau chiar folosirea altrei metode de învățare (precum KNN).